

→ Multi-head Attention

$$X \times W^Q = Q$$

$$X \times W^K = K$$

$$X \times W^V = V$$

$$\text{Softmax}\left(\frac{Q \times K^T}{\sqrt{d_k}}\right) V = Z$$

↳ Attention head

Attention head 0

Q_0 W_0^Q

W_0^K W_0^K

V_0 W_0^V

Z_0

Attention head 1

Q_1 W_1^Q

K_1 W_1^K

V_1 W_1^V

Z_1

Attention head N

Q_n W_n^Q

K_n W_n^K

V_n W_n^V

Z_n